

Big Data Analysis

Big Data

- ① A lot of data
- data Complexity
- ② Variety of data
- ③ data Streaming

- 1 → Structured
- 4 → Unstructured XML
- 2 → Semistructured
- 3 → Quasi Structured

web search history  
web click stream

test3 V model:

Volume Variety Velocity

Value of data  $\equiv$  Business Value4 V model:

Volume Variety Velocity Value

Certainty of data  $\equiv$  Veracity5 V model:

Volume Variety Velocity Value Veracity

R lang.  
 $\equiv$   
Tooltest  
ch/c of Big dataModule 4

distribution of data →

time lines → rate of data creation, every day/month/year/sec

data Analy scientist

→ must be doubtful

→ must be aware of lots of fields

OLAP → Online Analytical Processing

test  
slide 17Data Island → not connected, easy, any one can make it  
not secured, replication may exist  
(redundancy), data can get lostdata warehouse → secure, Backup exists, connected  
DB Administrator, is in control, No replication  
→ responsible for permissions



Sandbox → Analyst owns it not Administrator  
 Big Size, not small  
 Variety of data, more accurate results.  
 Shadow like sys. → I see the data in the Data warehouse but don't have a replication of it & can't change it without permission

### mini-Case Study

requirements

- more security
- relations between ~~connected~~ currencies
- laws / legal requirements
- distributed storage
- controlled access to data from all Banks
- more services (ATM) (Employees)

requirements from Analyst point of view:

- Velocity
- Variety (critical stuff)
- Volume
- Tools & methods to manage the data
- data driven (manages based on the data he gets)

KPI ≡ Key Performance Indicator

TP True +ve  
 FN False -ve  
 TN True -ve  
 FP False +ve

$$\frac{TP + TN}{\text{total}} = \text{accuracy}$$

Business intelligence, data Analyst // oral test

- Skills  
 - how can they help you

Business

Business Intelligence → Just an indicator  
 doesn't take decisions  
 uses KPI, deals with structured data

Data Analyst → ~~gives~~ Analysis &  
 gives decisions

High knowledge in AI, machine learning, Math, data mining